

---

# Melody Extraction Project Proposal

---

**Justin Mao-Jones**  
Center for Data Science  
New York University  
justinmaojones@nyu.edu

**Junbo Zhao**  
Center for Data Science  
New York University  
j.zhao@nyu.edu

**Rita Li**  
Center for Data Science  
New York University  
ml4713@nyu.edu

## 1 Introduction

We propose to build and evaluate machine learning models for the task of extracting melody from digital music files. Over the past decade, melody extraction has been an active research area in the music information retrieval research community [1]. It has many applications, including query by humming (e.g. hum a song into your phone and an app tells you what song it is), cover song identification, genre classification, mood classification, etc.

Generally speaking, melody is the predominant pitch in a piece of music that captures the essence of a song. As a motivating example, the melody is the tune one might hum when asked "what does the song sound like?" Unfortunately, there does not seem to be a precise definition of melody. As a machine learning task, we require some sort of precise definition, and for this project we adopt the definition described by [1]:

Melody is the fundamental frequency<sup>1</sup> from musical content with a lead voice or instrument. Melody extraction is the estimation of this melody from a single source.

While this definition is still open to interpretation, it can be used by human experts to generate the melody of a piece of music. Note that we constrain the definition of melody extraction to a single source, meaning that the melody is only coming from a single lead voice or instrument. The intuition behind this simplification is that it can make the task of melody extraction easier.

Melody extraction requires a set of labeled data. For this project, we propose to use MedleyDB [2], a dataset of annotated, royalty-free multitrack recordings that was curated primarily to support research on melody extraction. In section X, we discuss the dataset further.

Many approaches to melody extraction have been attempted [1], including pure signal processing [3], dynamic programming [4], support vector machines [5], and hidden markov models [6] [8]. In this project, we plan to utilize SVM, HMM, convolutional neural networks, and sparse coding. In section Y, we discuss these approaches further.

---

<sup>1</sup>Any audio signal can be represented as a sum of a series of sinusoids. The fundamental frequency of a signal defined as the lowest frequency in the series. It can be derived through the Fourier Transform.

## 1.1 Data

MedleyDB [2] consists of 122 multitracks, including stereo quality mixed audio, melody annotations, and stems<sup>2</sup>. The multitracks include songs from a variety of genres, including Singer/Songwriter, Classical, Rock, World/Folk, Fusion, Jazz, Pop, Musical Theatre, Rap.

The audio files are provided in WAV format (44.1 kHz, 16 bit). In other words, each audio file contains 44,100 digital audio samples per second. Each audio file is accompanied by a single source melody annotation, provided in csv format. A melody label is a number that represents the predominant frequency over a pre-defined window of time. There are 256 melody labels per second. Each melody label overlaps roughly 172 audio samples.

One of the benefits of MedleyDB is that it was carefully curated to provide a complete set of melody annotations combined with high-quality songs. Thus, we are operating under the assumption that we do not have missing data and that the melody annotations perfectly overlap with the audio samples. Unfortunately, it would not be practically feasible to check this assumption thoroughly, but we can do simple checks such as check the length of song against the length of its corresponding melody annotation.

An open question for our project team is whether or not we have enough data. On the one hand, 122 songs does not sound large. On the other hand, 256 melody labels per second, at an average of 3 minutes per song, corresponds to over 5 million training samples. So the question is whether or not there is enough variety in this dataset to generalize well to other datasets. Fortunately, there are additional datasets that we can utilize, such as RWC [7]. Time willing, we will try to add these datasets to our project.

## 2 Problem Definition

As a machine learning task, melody extraction has two components:

1. Voicing detection (i.e. classification of whether or not the melody is present),
2. Melody pitch tracking.

The first task derives from the fact that sometimes there is no melody. The tasks can be approached in separate algorithms or in the same algorithm. Melody pitch tracking seems to be the more difficult task of the two.

### 2.1 Melody Pitch Tracking

Melody pitch tracking is the task of identifying the predominant frequency during a *frame* (a small time interval). We are effectively constrained to the same frames used in the MedleyDB melody annotations, and thus our frames are roughly 4ms.

An important modeling question is whether to predict melody labels through regression or classification. Regression can seem to be a natural fit, because frequency lies on a continuous spectrum. However, human music tends to be composed on a discrete scale, i.e. musical notes. We believe that regression could be used to predict melodies, and predictions could be refined by "rounding" the predictions to the nearest note.

If using classification, then we would need to model all possible notes. For example, we could use the 88 keys found on a standard piano as the domain of our labels. Typically, most songs may not even cover this entire spectrum, and so we could reduce the size of our classification to only those notes present in our data.

---

<sup>2</sup>In a recording session, there is a separate microphone for each instrument (or sets of instruments), and thus there are separate recordings. For example, the singer is recorded separately from the guitar. A stem is one of these separate recordings. When added together the mix will sound like a complete song.

### 3 Evaluation Metrics For Model Performance

Given that melody extraction consists of two tasks, it is natural to evaluate each task separately and together, thus yielding three different evaluation metrics. Typically, researchers use accuracy to evaluate performance<sup>3</sup>, and so we will do the same here:

1. Voicing detection accuracy - TPR, i.e. probability that a frame that is actually voiced is predicted to be voiced,
2. Raw pitch accuracy - probability of a correct pitch value (within 1/4 tone<sup>4</sup>) given that the frame is voiced,
3. Overall accuracy =  $\frac{TPC+TN}{TO}$   
TP = total number of frames correctly predicted as voiced  
TPC = total number of TP frames in which pitch was also correctly predicted  
TN = total number of frames correctly predicted as unvoiced  
TO = total number of predictions

## 4 Approaches

### 4.1 Feature Extraction

Since Music Information Retrieval (MIR) is closely related to the Speech Recognition (SR) community, it naturally absorbs feature descriptors and machine learning techniques from the Speech community, which has a relatively larger amount of literature. We briefly studied the marriage between the communities and settled a few methods that might be appropriate.

One of the difficulties in both MIR and SR is that, in order to progress in model performance, task-specific handcrafted features are sometimes needed. We acknowledge this, but have decided that we will start with off-the-shell features such as Mel-Frequency Cepstral Coefficients (MFCC) [1], Short-time Fourier Transform (STFT) and multi-resolution FFT (MRFFT) [8].

- **STFT.** STFT is a widely used signal preprocessing technique. In general, the signal is chunked into frames where STFT is applied to each chunked frame, with a window length typically assigned as 50 and 100ms. Frames can overlap in a sliding window fashion.
- **MRFFT.** Resolution issues inherently arise with Fourier transform. MRFFT overcomes this by taking frequency spectrum out of multi-resolution windows.
- **MFCC.** This has been the dominant feature descriptor in Speech community over the past 30 years. It is basically a linear cosine encoding of the log power spectrum on a mel-scale of frequency which biologically originated from human's ears.
- **Dictionary Learning.** Dictionary learning is an adaptive content-based feature self-learning approach. Its goal basically is to get local descriptions by learning a linear combination of a pre-defined dictionary. The weights of the words in the dictionary are the new representation of the local window. The process of getting the dictionary is unsupervised; K-means and K-SVD [9] are two common methods to obtain this dictionary.

### 4.2 Classifier or Regressor

A natural next step is to fit extracted features into a classifier or regressor that employs the power of Machine Learning. Hidden Markov Model (HMM) and Support Vector Machine (SVM) are our primary elections.

- **SVM.** Both linear kernel and RBF kernel would be put into uses.
- **HMM.** HMM is a directed model in which the system being modeled is regarded as Markov process with hidden states (unobserved).

<sup>3</sup>[http://www.music-ir.org/mirex/wiki/2014:Audio\\_Melody\\_Extraction](http://www.music-ir.org/mirex/wiki/2014:Audio_Melody_Extraction)

<sup>4</sup>a tone is a step between two keys on a piano; e.g. G is a tone higher than F.

### 4.3 End-to-end learning

Inspired by the recent success of deep learning in SR, an alternative approach is to apply end-to-end learning techniques which derives useful information from the raw audio waves [J4]. Convolution Neural Network (CNN) and Recurrent Neural Network (RNN) are two predominant choices modeling sequential data. Both of them can be applied to the raw audio input or Fourier transformed features and trained by gradient based method using back propagation.

- **CNN.** CNN operates between Temporal Convolution, Temporal Pooling and Non-linear activations. On top of the network framework, a classifier or regressor would be placed to finish the job of classifying or regressing.
- **RNN.** Differing from feedforward neural network, RNN uses the internal memory to process signals in such a way that RNN output from previous time steps are used as inputs for subsequent time steps.

### References

- [1] J. Salamon, E. Gomez, D. P. W. Ellis and G. Richard, "Melody Extraction from Polyphonic Music Signals: Approaches, Applications and Challenges", IEEE Signal Processing Magazine, 31(2):118-134, Mar. 2014.
- [2] R. Bittner, J. Salamon, M. Tierney, M. Mauch, C. Cannam and J. P. Bello, "MedleyDB: A Multitrack Dataset for Annotation-Intensive MIR Research", in 15th International Society for Music Information Retrieval Conference, Taipei, Taiwan, Oct. 2014.
- [3] J. Salamon and E. Gmez, "Melody Extraction from Polyphonic Music Signals using Pitch Contour Characteristics", IEEE Transactions on Audio, Speech and Language Processing, 20(6):1759-1770, Aug. 2012.
- [4] V. Rao and P. Rao, Vocal melody extraction in the presence of pitched accompaniment in polyphonic music, IEEE Trans. Audio, Speech, Lang. Processing, vol. 18, no. 8, pp. 2145-2154, Nov. 2010.
- [5] G. Poliner and D. Ellis, A classification approach to melody transcription, in Proc. 6th Int. Conf. Music Information Retrieval, London, Sept. 2005, pp. 161-166.
- [6] M. Ryynnen and A. Klapuri, Automatic transcription of melody, bass line, and chords in polyphonic music, Comput. Music J., vol. 32, no. 3, pp. 72-86, 2008.
- [7] Masataka Goto, Hiroki Hashiguchi, Takuichi Nishimura, and Ryuichi Oka: RWC Music Database: Popular, Classical, and Jazz Music Databases, Proceedings of the 3rd International Conference on Music Information Retrieval (ISMIR 2002), pp.287-288, October 2002.
- [8] T.-C. Yeh, M.-J. Wu, J.-S. Jang, W.-L. Chang, and I.-B. Liao, A hybrid approach to singing pitch extraction based on trend estimation and hidden Markov models, in IEEE Int. Conf. Acoustics, Speech, and Signal Processing (ICASSP), Kyoto, Japan, Mar. 2012, pp. 457-460.
- [9] Aharon, Michal and Elad, Michael and Bruckstein, Alfred, "K-SVD: An Algorithm for Designing Overcomplete Dictionaries for Sparse Representation," in IEEE Transactions on Signal Processing 2006, Vol 54 No 11 pp. 4311.